

## Multiple-bias modelling for analysis of observational data

Sander Greenland

*University of California, Los Angeles, USA*

[*Read before The Royal Statistical Society on Wednesday, September 29th, 2004, the President, Professor A. P. Grieve, in the Chair*]

**Summary.** Conventional analytic results do not reflect any source of uncertainty other than random error, and as a result readers must rely on informal judgments regarding the effect of possible biases. When standard errors are small these judgments often fail to capture sources of uncertainty and their interactions adequately. Multiple-bias models provide alternatives that allow one systematically to integrate major sources of uncertainty, and thus to provide better input to research planning and policy analysis. Typically, the bias parameters in the model are not identified by the analysis data and so the results depend completely on priors for those parameters. A Bayesian analysis is then natural, but several alternatives based on sensitivity analysis have appeared in the risk assessment and epidemiologic literature. Under some circumstances these methods approximate a Bayesian analysis and can be modified to do so even better. These points are illustrated with a pooled analysis of case-control studies of residential magnetic field exposure and childhood leukaemia, which highlights the diminishing value of conventional studies conducted after the early 1990s. It is argued that multiple-bias modelling should become part of the core training of anyone who will be entrusted with the analysis of observational data, and should become standard procedure when random error is not the only important source of uncertainty (as in meta-analysis and pooled analysis).

**Keywords:** Bayesian statistics; Confidence profile method; Confounding; Epidemiologic methods; Leukaemia; Magnetic fields; Meta-analysis; Meta-statistics; Monte Carlo methods; Observational data; Odds ratio; Relative risk; Risk analysis; Risk assessment; Sensitivity analysis

### 1. Introduction

#### 1.1. *The problem*

In their discussion of observational data analysis, Mosteller and Tukey (1977), page 328, said standard errors ‘cannot be expected to show us the indeterminacies and uncertainties we face’. More recently, a prize winning paper by Maclure and Schneeweiss (2001) described how random error is but one component in a long sequence of distortive forces leading to epidemiologic observations and is often not the most important. Yet conventional analyses of observational data in the health sciences (as reviewed, for example, in Rothman and Greenland (1998), chapters 12–17) can be characterized by a two-step process that quantifies only random error—

- (a) employ frequentist statistical methods based on the following assumptions, which may be grossly violated in the application but are not testable with the data under analysis:

*Address for correspondence:* Sander Greenland, Departments of Epidemiology and Statistics, University of California at Los Angeles, 22333 Swenson Drive, Topanga, CA 90290, USA.  
E-mail: lesdomes@ucla.edu

- (i) the study exposure is randomized within levels of controlled covariates (sometimes replaced by a practically equivalent assumption of 'no unmeasured confounders' or 'ignorability of treatment assignment');
  - (ii) selection, participation and missing data are random within levels of controlled covariates;
  - (iii) there is no measurement error (occasionally, an unrealistically restrictive error model is used to make a correction, which can do more harm than good; see Wacholder *et al.* (1993));
- (b) address possible violations of assumptions (i)–(iii) with speculative discussions of how each might have biased the statistical results. If they like the statistical results from the first step, researchers will argue that these biases are inconsequential, rarely offering evidence to that effect (Jurek *et al.*, 2004). However if they dislike their results they may focus on possible biases and may even write whole articles about them (e.g. Hatch *et al.* (2000)).

In practice, the second step is often skipped or fails to address more than one or two assumptions (Jurek *et al.*, 2004). The assumptions in the first step can be replaced by the slightly weaker assumption that any biases from violations of (i)–(iii) cancel, but appeal to such cancellation seems wishful thinking at best.

Paul Meier (personal communication) and others have defended conventional results (derived under step (a)) as 'best case' scenarios that show the absolute minimum degree of uncertainty that we should have after analysing the data. Unfortunately, the above assumptions are far too optimistic, in that they produce misleadingly narrow interval estimates precisely when caution is most needed (e.g. in meta-analyses and similar endeavours with potentially large policy impact, as illustrated below). Worse, the illusory precision of conventional results is rarely addressed by more than intuitive judgments based on flawed heuristics; see Section 4.3.

Another defence is that conventional results merely quantify random error. This defence overlooks the fact that such quantification is hypothetical and hence questionable when no random sampling or randomization has been employed and no natural random mechanism has been documented. Conventional (frequentist) statistics are still often touted as 'objective', even though in observational epidemiology and social science they rarely meet any criterion for objectivity (such as derivation from a mechanism that is known to be operative in the study). This belief has resulted in an unhealthy obsession with random error in both statistical theory and practice. A prime example, which is often lamented but still very much a problem, is the special focus that most researchers give to 'statistical significance'—a phrase whose very meaning in observational studies is unclear, owing to the lack of justification for conventional sampling distributions when random sampling and randomization are absent.

The present paper is about the formalization of the second step to free inferences from dependence on the highly implausible assumptions that are used in the first step and the often misleading intuitions that guide the second step. Although I limit the discussion to observational studies, the bias problems that I discuss often if not usually arise in clinical trials, especially when non-compliance or losses occur, and the methods described below can be brought to bear on those problems.

### 1.2. An overview of solutions

An assessment of uncertainty due to questionable assumptions (uncertainty analysis) is an essential part of inference. Formal assessments require a model with parameters that measure departures from those assumptions. These parameters govern the bias in methods that rely

on the original assumptions; hence I shall call the parameters *bias parameters*, the model for departures a *bias model* and departures from a particular assumption a *bias source*.

Statistical literature on bias models remains fragmented; most of it deals with just one bias source, and the bias model is often used only for a sensitivity analysis (which displays bias as a function of the model parameters), although occasionally it becomes part of a Bayesian analysis. In contrast, the literature on risk assessment and decision analysis has focused on accounting for all major sources of uncertainty (Morgan and Henrion, 1990; Crouch *et al.*, 1997; Vose, 2000; Draper *et al.*, 2000). Most notable in the health sciences are the confidence profile method (Eddy *et al.*, 1992), which incorporates bias models into the likelihood function, analyses based on non-ignorable non-response models with unknown bias parameters (Little and Rubin, 2002), and Monte Carlo sensitivity analysis (MCSA), which samples bias parameters and then inverts the bias model to provide a distribution of 'bias-corrected' estimates (Lash and Silliman, 2000; Powell *et al.*, 2001; Lash and Fink, 2003; Phillips, 2003; Greenland, 2003a, 2004a; Steenland and Greenland, 2004).

### 1.3. Outline of paper

The next section gives some general theory for bias modelling that encompasses frequentist (sensitivity analysis), Bayesian and MCSA approaches. The theory gives a formal perspective on MCSA and suggests ways to bring it closer to posterior sampling. In particular, it operationalizes the sequential bias factor approach (Maclure and Schneeweiss, 2001) in a form that approximates Gibbs sampling under certain conditions and explains the similarity of Bayesian and MCSA results that are seen in published examples (Greenland, 2003a; Steenland and Greenland, 2004). Section 3 analyses 14 studies of magnetic fields and childhood leukaemia, extending a previous analysis (Greenland, 2003a) by adding new data, providing more detail in illustration and extending the bias model to include classification error. Classification error is a large source of uncertainty due to an absence of data on which to base priors, and due to the extreme sensitivity of results over reasonable ranges for the bias parameters. Section 4 discusses some problems in interpreting and objections to bias modelling exercises; it argues that many of the criticisms apply with even more force to conventional analyses, and that the status of the latter as expected and standard practice in observational research is unwarranted. That section can be read without covering Sections 2 and 3, and I encourage readers who are uninterested in details to skim those two sections and to focus on Section 4.

## 2. Some theory for observational statistics

### 2.1. Model expansion to include bias parameters

To review formal approaches to the bias problem, suppose that the objective is to make inferences on a target parameter  $\theta = \theta(\alpha)$  of a population distribution parameterized by  $\alpha$ , using an observed data array  $A$ . Conventional inference employs a model  $L(A; \alpha)$  for the distribution of  $A$  given  $\alpha$  and some background assumptions that are sufficient to identify  $\theta$  from  $A$ , such as 'randomization of units to treatment', 'random sampling of observed units and of data on those units' and 'no measurement error' (step (a)(i)–(a)(iii) above). Most statistical methodology concerns extensions of basic models, tests and estimators to complex sampling, allocation and measurement structures. The identification of  $\theta$  is retained by treating these structures as known or as jointly identifiable with  $\theta$  from  $A$  under the assumed model, making assumptions as necessary to ensure identifiability (e.g. assumptions of 'no unmeasured confounders', 'missing at random' and 'ignorable non-response').

On the basis of assumptions about their operation, the effects of bias sources on  $L$  may be modelled by using a bias parameter vector  $\eta$ , so that the data distribution is represented by an expanded model  $L(A; \alpha, \eta)$ . Examples include most models for uncontrolled confounding and response bias (such as non-ignorable treatment assignment models and non-response models with unknown parameters) (e.g. Leamer (1974, 1978), Rubin (1983), Copas and Li (1997), Robins *et al.* (1999), Gelman *et al.* (2003), Little and Rubin (2002), Rosenbaum (2002) and Greenland (2003a)). Typically,  $\alpha$  is not even partially identified without information on  $\eta$ , in that every distinct distribution in the family  $L(A; \alpha, \eta)$  can be generated from a given  $\alpha$  by finding a suitable  $\eta$ , and a prior that is uniform in  $\eta$  leads to  $p(\alpha|A) \approx p(\alpha)$ . Thus, inferences about  $\alpha$  are infinitely sensitive to  $\eta$ , and  $L(A; \alpha, \eta)$  is uninformative for  $\alpha$  (and hence  $\theta$ ) without prior information on  $\eta$ . In the same manner,  $\eta$  is not identified without information on  $\alpha$ , so that a prior that is uniform in  $\alpha$  leads to  $p(\eta|A) \approx p(\eta)$ .

I shall consider only large sample behaviour. 'Unbiased' will be mean  $\sqrt{n}$ -consistent uniformly asymptotically unbiased normal, as is customary in much epidemiologic statistics. For simplicity I shall assume that  $\alpha$  fully specifies the population distribution, but the theory can be extended to semiparametric models by using familiar modifications of likelihood (conditional, partial, etc.). I shall also assume that any necessary regularity conditions hold, e.g. the joint parameter space of  $(\alpha, \eta)$  is the product of the marginal spaces of  $\alpha$  and  $\eta$ , and all models and functions are smooth in their arguments and parameters. Given these conditions, conventional estimators of  $\theta$  extend naturally to the expanded model. For example, suppose that  $\hat{\theta}_\eta$  and  $\hat{s}_\eta$  are the maximum likelihood estimator (MLE) of  $\theta$  and its estimated standard error obtained from  $L(A; \alpha, \eta)$  when the bias parameter is fixed at a known value  $\eta$ . Under the models that are used here,  $\hat{\theta}_\eta$  is unbiased for  $\theta$  and  $\hat{\theta}_\eta \pm 1.96\hat{s}_\eta$  is a large sample 95% confidence interval when the model and value of  $\eta$  that are used in it are correct. Parameterizing  $L$  so that  $\eta=0$  corresponds to no bias,  $L(A; \alpha, 0)$  then represents the conventional analysis distribution,  $\hat{\theta}_0$  is the conventional estimator,  $E(\hat{\theta}_0) - \theta$  is its (asymptotic) bias and  $\hat{\theta}_0 - \hat{\theta}_\eta$  is its estimated bias given  $\eta$ .

## 2.2. Sensitivity analysis

Because  $\eta$  is unknown and not identified, the preceding results are of little use by themselves. Sensitivity analyses display how statistics like  $\hat{\theta}_\eta$  and derived confidence limits and  $P$ -values vary with  $\eta$ . Epidemiologic examples date back at least to Cornfield *et al.* (1959), and since then the methods have been extended to many settings (e.g. Eddy *et al.* (1992), Greenland (1996), Copas and Li (1997), Copas (1999), Robins *et al.* (1999), Little and Rubin (2002) and Rosenbaum (2002)). Yet sensitivity analysis remains uncommon in health and medical research reports. This is not surprising, given the lack of motivation for its use and its relative unfamiliarity: sensitivity analysis is mentioned in few journal instructions or statistics text-books. As with informal discussions, those sensitivity analyses that are published rarely examine more than one bias at a time and so overlook interactions, such as those that arise from covariate effects on classification errors (Greenland and Robins, 1985; Flegal *et al.*, 1991; Lash and Silliman, 2000).

To address this concern we can use a model  $L(A; \alpha, \eta)$  that incorporates multiple bias sources; indeed, my thesis is that realistic uncertainty analyses of observational data must do so. Nonetheless, the difficulty of examining a grid beyond three dimensions necessitates some sort of summarization over the sensitivity results. If (as here) the net bias in the conventional estimator  $\hat{\theta}_0$  is not constrained by the data, any reasonable summary will be determined entirely by the choices of values for  $\eta$  and so will be arbitrarily sensitive to those choices (Greenland, 1998). Furthermore, apparent data constraints on bias and hence on inference can depend entirely on the structure of the data model and can disappear after only minor elaboration (Poole and Greenland, 1997).

These problems lead to another obstacle for the adoption of sensitivity analysis: its potential for arbitrary or nihilistic output. In the present setting, for any preselected value  $\theta_v$  for  $\theta$ , we can find a value for  $\eta$  that yields  $\hat{\theta}_\eta = \theta_v$ ; thus, any output pattern can be produced by manipulating  $\eta$ . Although an arbitrary or purely manipulative analysis (one that displays a pattern that is preselected by the analyst) might be obvious in a simple case, it might not be so obvious with multiple bias sources.

To summarize: sensitivity analysis only describes the dependence of statistics on  $\eta$ .  $\eta$  is often of high dimension. The complexity of the dependence can render sensitivity analyses difficult to present without drastic (and potentially misleading) simplifications. Furthermore, sensitivity analysis may exclude no possible value for  $\theta$ : results can be infinitely sensitive to  $\eta$ , and hence without some constraint on  $\eta$  the analysis will only display this fact. The constraints chosen can play a pivotal role in the appearance of the results, and informed choices essentially correspond to a prior for  $\eta$  (Greenland, 1998, 2001a).

### 2.3. Bayesian analysis and Monte Carlo sensitivity analysis

One way to address the limits of sensitivity analysis is to specify explicitly a prior density  $p(\alpha, \eta)$  and base inferences for  $\theta = \theta(\alpha)$  on the marginal posterior

$$p(\alpha|A) \propto \int L(A; \alpha, \eta) p(\alpha, \eta) d\eta$$

(Leamer, 1974, 1978; Eddy *et al.*, 1992; Graham, 2000; Little and Rubin, 2002; Gustafson, 2003; Greenland, 2001a, 2003a). To account for shared prior information (and the resulting prior correlations) between components of  $\eta$ , the bias parameter  $\eta$  may itself be modelled as a function of known covariates and unknown hyperparameters  $\beta$ , resulting in a hierarchical bias model (Greenland, 2003a), as below. None-the-less, many health researchers reject formal Bayesian methods as too difficult if not philosophically objectionable: analytic solutions for  $p(\theta|A)$  involving just one bias source can appear formidable (Eddy *et al.*, 1992; Graham, 2000; Gustafson, 2003), and sampler convergence remains crucial yet extremely technical (Gelman *et al.*, 2003; Gustafson, 2003).

An easier alternative specifies only a marginal prior  $p(\eta)$  for the bias parameters, samples  $\eta$  from this prior, computes  $\hat{\theta}_\eta$  and  $\hat{s}_\eta$  from each sample and summarizes the resulting distribution of  $\hat{\theta}_\eta$  and of statistics derived from  $\hat{\theta}_\eta$  and  $\hat{s}_\eta$ . The  $\hat{\theta}_\eta$  that are generated by this MCSA have various uses. The distribution of  $\hat{\theta}_0 - \hat{\theta}_\eta$  estimates the distribution of net bias under  $p(\eta)$ , and the distribution of  $\hat{\theta}_\eta$  can be compared with the sampling distribution of  $\hat{\theta}_0$  to measure the relative importance of bias uncertainty and random error. Standard errors shrink as data accumulate and hence bias uncertainty grows in importance and eventually dominates uncertainty due to random error. As will be illustrated, the comparison of bias uncertainty and random error can reveal that the benefits of study replication diminish far below those indicated by conventional power calculations, for the latter ignore bias uncertainty.

With modification, MCSA can also provide approximate posterior inferences. Suppose that, for each  $\eta$ ,  $\hat{\theta}_\eta$  is approximately efficient (e.g. is the MLE),  $p(\alpha|\eta)$  is approximately uniform and  $p(\eta|A) \approx p(\eta)$ . We then have approximately  $p(\alpha, \eta) \propto p(\eta)$  and

$$p(\alpha|A, \eta) \propto L(A; \alpha, \eta) p(\alpha, \eta) / p(\eta|A) \propto L(A; \alpha, \eta),$$

and

$$p(\theta|A, \eta) \propto \int_{\theta(\alpha)=\theta} L(A; \alpha, \eta) d\alpha$$

with the latter approximately normal( $\hat{\theta}_\eta, \hat{s}_\eta^2$ ) (Gelman *et al.* (2003), chapter 4). Thus, the MCSA procedure can be modified to approximate sampling from

$$p(\theta|A) = \int p(\theta|A, \eta) p(\eta|A) d\eta$$

by

- (a) drawing  $\eta$  from  $p(\eta)$
- (b) computing  $\hat{\theta}_\eta$  and  $\hat{s}_\eta^2$ , and
- (c) redrawing  $\hat{\theta}_\eta$  from a normal( $\hat{\theta}_\eta, \hat{s}_\eta^2$ ) distribution or (equivalently) adding a normal( $0, \hat{s}_\eta^2$ ) disturbance to  $\hat{\theta}_\eta$  (Greenland, 2003a).

If  $\eta$  partitions into  $\eta_k$  that are estimable given  $\alpha$  and the remaining components  $\eta_{-k}$ , the algorithm generalizes to arbitrary  $p(\alpha, \eta)$  by cycling among  $p(\alpha|A, \eta)$  and the  $p(\eta_k|A, \alpha, \eta_{-k})$ , drawing from an approximate normal distribution at each step, whence it can be seen as a large sample approximation to Gibbs sampling.

To avoid normal approximations, some researchers resample the data as well as  $\eta$  at each trial (Lash and Fink, 2003). Naïve resampling (i.e. bootstrapping from the empirical data distribution) does, however, have its own small sample problems (Efron and Tibshirani, 1993); for example, in tabular data it leaves empty observed cells as 0s and so will never visit some points in the support of the sampling distribution. To remove these 0s, we may resample the data from a smoothed table, then smooth each resample with the same procedure as that used to smooth the original data.

#### 2.4. Some useful specializations for discrete data

Suppose now that  $A$  represents a count vector for a multiway cross-classification of the data. Conventional approaches model  $A$  with a distribution  $L\{A; E(A; \alpha)\}$  that depends on the population parameters only through the expected counts  $E(A; \alpha)$ . Suppressing  $\alpha$  in the notation, one extension takes  $E_\eta \equiv E(A; \alpha, \eta)$ , with  $E_0 = E(A; \alpha, 0)$  the counts expected in the absence of bias, so that the expanded model can be written  $L(A; E_\eta)$ . Note that  $E_\eta$  is an estimable quantity even though  $\eta$  and  $E_0$  are not separately identified. For example, with no constraint on  $\eta$  or  $E_0$ , the multinomial MLE of  $E_\eta$  is the observed  $A$ . Hence we can model  $E_\eta$  directly, as will be done below for smoothing purposes.

For some models,  $E_\eta = G_\eta(E_0)$  where the ‘bias function’  $G_\eta$  is a family of mappings indexed by  $\eta$  and  $G_0$  is the identity;  $\hat{\theta}_\eta$  may then be taken as the MLE of  $\theta$  from  $L\{A; G_\eta(E_0)\}$ . These models can be especially simple. For example, if the only bias source is non-response and  $\eta_R$  is the vector of log-response-rates within cells of the observed cross-classification,  $E_\eta = G_\eta(E_0; \eta_R) = \text{diag}\{\exp(\eta_R)\}E_0$ ;  $\eta_R$  thus becomes a log-linear offset to the conventional model for  $E_0$ , and  $\hat{\theta}_\eta$  is the MLE from the offset model for  $E_\eta$ . Confounding can also be represented by a log-linear offset  $\eta_C$  which can be added to the response-bias offset, although this offset is a non-linear function of unmeasured covariate effects (see the example below); in that case  $E_\eta = G_\eta(E_0; \eta) = B_\eta E_0$  where  $B_\eta = \text{diag}\{\exp(\eta_C + \eta_R)\}$  (Greenland, 2003a).

Next, suppose that  $q_{ij}$  is the probability that a unit will be classified in cell  $i$  of  $A$  given that it should be in cell  $j$ . With  $Q$  the matrix of  $q_{ij}$  and  $\eta_M = \text{vec}(Q)$ , one model for confounding and non-response followed by misclassification would be  $E_\eta = G_\eta(E_0; \eta) = B_\eta E_0$  where  $B_\eta = Q \text{diag}\{\exp(\eta_C + \eta_R)\}$ . Alternatively, suppose that  $p_{ij}$  is the probability that a unit should be in cell  $i$  given that it is classified in cell  $j$ ; with  $P$  the matrix of  $p_{ij}$  and  $\eta_M = \text{vec}(P)$  we have  $E_\eta = G_\eta(E_0; \eta) = B_\eta E_0$  where  $B_\eta = P^{-1} \text{diag}\{\exp(\eta_C + \eta_R)\}$ .

Without validation data identifying  $P$ , acceptable assumptions or priors about misclassification more often concern  $Q$  than  $P$ , as below. An important difference between the  $Q$ - and  $P$ -models is that  $A$  is informative for  $Q$  even without information about  $\alpha$ ; hence use of  $Q$  may lead to  $p(\eta|A) \neq p(\eta)$ . For example, a non-zero observed cell  $i$  implies that  $q_{ij} > 0$  for some  $j$ ; further assumptions can lead to stronger constraints on  $Q$ . In contrast,  $A$  alone does not constrain  $P$ . Thus, the above arguments for MCSA as an approximation to posterior sampling do not strictly apply under the  $Q$ -model unless the support of the prior  $p(Q)$  falls within the identified bounds.

2.5. Sequential correction

If  $G_\eta$  is invertible, a ‘bias-corrected’ estimator  $\hat{\theta}_{0\eta}$  of  $\theta$  can be obtained by applying a conventional estimator  $\hat{\theta}_0$  to the ‘corrected data’  $F(A; \eta) \equiv G_\eta^{-1}(A)$ , where  $F(A; 0) = A$  (Lash and Fink, 2003).  $\hat{\theta}_{0\eta}$  is an unbiased estimator of  $\theta$  given that  $\eta$  and the model are correct, but the ‘standard error’  $s_{0\eta}$  for  $\hat{\theta}_{0\eta}$  that is obtained by applying a conventional variance estimator to  $F(A; \eta)$  is not generally valid (see below).

$F(A; \eta)$  is typically derived from separate correction formulae in conventional sensitivity analyses (Rothman and Greenland (1998), chapter 19). There are many formulae  $F_C(\cdot; \eta_C)$ ,  $F_R(\cdot; \eta_R)$  and  $F_M(\cdot; \eta_M)$  that correct for confounding, response bias and misclassification. For example, with  $\eta_R$  the vector of log-response-rates, a correction that adjusts the relative frequencies for non-response is  $F_R(A; \eta_R) = \text{diag}\{\exp(c - \eta_R)\}A$ , where  $c$  is a log-normalizing-constant vector to preserve  $A$ -margins that are fixed by design. With  $P$  and  $Q$  as above, correction formulae for misclassification include  $F_M(A; \eta_M) = PA$  and  $F_M(A; \eta_M) = Q^{-1}A$ ; these formulae automatically preserve fixed margins if (as is often the case) misclassification can only occur within the strata that are defined by those margins, for then  $p_{ij} = q_{ij} = 0$  when  $i$  and  $j$  are in different strata, and hence  $P$  and  $Q$  are block diagonal when the indices are ordered by stratum.

Use of the observed counts in the formulae corresponds to using  $E_\eta = A$ , which is a saturated model for  $E_\eta$ . Some formulae (like that based on  $Q$ ) can yield impossible (e.g. negative) corrected counts for certain values of  $\eta$ , especially if 0s are present in  $A$ , which lead to breakdown (division by 0) or wild behaviour of  $\hat{\theta}_{0\eta}$ . These problems can often be avoided by preliminary smoothing of  $A$  to remove non-structural 0s, e.g. by averaging  $A$  with a model-fitted count (which generalizes adding a constant to each cell (Bishop *et al.* (1975), chapter 12, and Good (1983), section 9.4), or by replacing  $A$  with a count that is expected under a nearly saturated model that preserves data patterns regardless of the statistical significance of the patterns (Greenland, 2004b).

Formulae can be applied in sequence to correct multiple biases, although the order of corrections is important if the formulae do not commute (Greenland, 1996). One can imagine each correction moving a step from the biased data back to the unbiased structure, as if hypothetically ‘unwrapping the truth from the data package’. For example, suppose that the data generation process is one in which causal effects (including effects of unmeasured confounders) generate population associations, subjects are sampled in a manner that is subject to non-response and finally the responding subjects are classified subject to error. This chronology suggests that we should correct misclassification first, then non-response, and then uncontrolled confounding. With  $\eta = (\eta_C, \eta_R, \eta_M)$ , the resulting bias-corrected counts  $F(A; \eta)$  are  $F_C[F_R\{F_M(A; \eta_M); \eta_R\}; \eta_C]$ . If the bias model is  $E_\eta = B_\eta E_0$ , we have  $F(A; \eta) = B_\eta^{-1}A$ ; for example, under the  $Q$ -model above,

$$B_\eta^{-1} = \text{diag}\{\exp(c - \eta_C - \eta_R)\}Q^{-1}.$$

The sequential correction approach can be simpler both conceptually and computationally than fully Bayesian or likelihood-based sensitivity approaches. We just plug the sampled bias

parameters ( $\eta_C, \eta_R, \eta_M$ ) into their respective formulae, apply the resulting corrections in proper sequence and compute  $\hat{\theta}_{0\eta}$  from the resulting  $F(A; \eta)$ , possibly replacing  $A$  by a smoothed count. If  $\hat{\theta}_0$  has a closed form (e.g. a Mantel–Haenszel estimator) then  $\hat{\theta}_{0\eta}$  will also be of closed form, resulting in a very rapid Monte Carlo procedure. In some examples such as that below, confounding and response bias corrections simplify to division of conventional stratum-specific odds ratio estimates by independent bias factors free of the data, leading to an even simpler and more rapid procedure. Finally, as with Bayesian analyses, in some simple cases the entire MCSA distribution has a closed form approximation (Greenland, 2001a).

**2.6. Sequential correction and posterior sampling**

The earlier arguments for MCSA as approximate posterior sampling hinge on the use of the MLE or an equivalent  $\hat{\theta}_\eta$  derived under the expanded model  $L(A; E_\eta)$  and so do not extend to the use of  $\hat{\theta}_{0\eta}$ . Suppose that under the conventional model  $\hat{\theta}_0$  is asymptotically equivalent to  $\theta(\hat{\alpha})$ , where  $\hat{\alpha}$  is an inverse variance weighted least squares estimator of  $\alpha$  from a regression of  $A$  on the classification axes, e.g. as when  $\alpha$  comprises log-linear model parameters and  $\hat{\theta}_0$  is the conventional MLE of a log-odds ratio. Then, if  $\eta=0$ ,  $\hat{\theta}_0$  is asymptotically efficient and first order equivalent to the conventional MLE. Because  $\hat{\theta}_{0\eta}$  treats  $F(A; \eta)$  as the observed counts, however, when  $\eta \neq 0$  the implicit weights are no longer the correct inverse variances and  $\hat{\theta}_{0\eta}$  is no longer efficient.

As an example, using maximum likelihood log-linear Poisson regression, the weight matrix for  $\ln\{F(A; \eta)\}$  which is implicit in  $\hat{\theta}_{0\eta}$  is  $W_{0\eta} = \text{diag}\{F(E_\eta; \eta)\}$ . The asymptotic inverse covariance matrix for  $\ln\{F(A; \eta)\}$  is, however,

$$W_\eta = \{D'_\eta \text{diag}(E_\eta) D_\eta\}^{-1}$$

where  $D_\eta = \partial[\ln\{F(E_\eta; \eta)\}]/\partial E_\eta$ . If  $F(A; \eta) = B_\eta^{-1} A$ , then  $D_\eta^{-1} = B_\eta \text{diag}(E_0)$  and hence  $W_\eta = B_\eta \text{diag}(E_0^2/E_\eta) B'_\eta \neq W_{0\eta} = \text{diag}(B_\eta^{-1} E_\eta) = \text{diag}(E_0)$  unless  $B_\eta$  is the identity. Furthermore, when  $B_\eta$  is diagonal (as when only confounding and response bias are corrected),  $W_\eta$  reduces to  $\text{diag}(E_\eta)$ , the weight matrix for the uncorrected regression, rather than to  $W_{0\eta}$ .

Because  $W_{0\eta}$  does simplify to  $W_\eta$  when  $\eta=0$ , the sequential estimator using  $\hat{\theta}_{0\eta}$  can be viewed as an approximation to the MLE  $\hat{\theta}_\eta$  in a neighbourhood of  $\eta=0$ . The Monte Carlo distribution of  $\hat{\theta}_{0\eta}$  over  $p(\eta)$  might thus be reasonably expected to approximate that of the MLE  $\hat{\theta}_\eta$  if  $p(\eta)$  is centred on zero and is not too dispersed. Alternatively, if  $\hat{\theta}_{0\eta}$  has an explicit weighted form, we could just estimate  $W_\eta$  directly and use that to compute  $\hat{\theta}_{0\eta}$ . Unfortunately, after misclassification correction,  $W_\eta$  is not diagonal, does not readily simplify along with the bias corrections and must be recomputed for each  $\eta$ . To avoid these problems, we could just use the diagonal matrix of uncorrected weights, which under the models that are used here would approximate the correct weights in a neighbourhood of  $\eta_M=0$  rather than just  $\eta=0$ . In examples based on the data below and with similar priors, the latter estimator augmented by a normal( $0, \hat{s}_0^2$ ) disturbance very closely approximated posterior distributions (Greenland, 2003a), so only this modified sequential approach will be illustrated.

**3. Magnetic fields and childhood leukaemia**

**3.1. The data**

The example data in Table 1 are taken from a pooled analysis of 12 pre-1999 case-control studies of residential magnetic fields and childhood leukaemia (Greenland, Sheppard, Kaune, Poole and Kelsh, 2000), plus two additional studies unpublished at the time that the analysis was done



**Table 1.** Summary data from 14 case-control studies of magnetic fields and childhood leukaemia

Reference	Country	Number of cases		Number of controls		Odds ratio (95% limits)
		>3 mG	Total	>3 mG	Total	
Coghill <i>et al.</i> (1996)	England	1	56	0	56	$\infty$
Dockerty <i>et al.</i> (1998)	New Zealand	3	87	0	82	$\infty$
Feychting and Ahlbom (1993)	Sweden†	6	38	22	554	4.53 (1.72,12.0)
Kabuto (2003)	Japan	11	312	13	603	1.66 (0.73,3.75)
Linet <i>et al.</i> (1997)	USA‡	42	638	28	620	1.49 (0.91,2.44)
London <i>et al.</i> (1991)	USA‡	17	162	10	143	1.56 (0.69,3.53)
McBride <i>et al.</i> (1999)	Canada‡	14	297	11	329	1.43 (0.64,3.20)
Michaelis <i>et al.</i> (1998)	Germany	6	176	6	414	2.40 (0.76,7.55)
Olsen (1993)	Denmark†	3	833	3	1666	2.00 (0.40,9.95)
Savitz <i>et al.</i> (1988)	USA‡	3	36	5	198	3.51 (0.80,15.4)
Tomenius (1986)	Sweden	3	153	9	698	1.53 (0.41,5.72)
Tynes and Haldorsen (1997)	Norway†	0	148	31	2004	0
UK Childhood Cancer Study Investigators (1999)	UK§	5	1057	3	1053	1.66 (0.40,6.98)
Verkasalo <i>et al.</i> (1993)	Finland†	1	32	5	320	2.03 (0.23,18.0)
Totals§§		115	4025	146	8740	1.69 (1.28,2.23)

†Calculated fields (the others are direct measurement).

‡120 V-60 Hz systems (the others are 220 V-50 Hz).

§Comparison of >4 mG versus  $\leq 2$  mG, excluding 16 cases and 20 controls at 2-4 mG.

§§The final column is the MLE of the common odds ratio (lower  $P=0.0001$ ; homogeneity  $P=0.24$ ).

(Kabuto, 2003; UK Childhood Cancer Study Investigators, 1999). Because the UK childhood cancer study did not supply individual data, its estimate compares the published categories of greater than 4 mG versus less than or equal to 2 mG. It is included here on the basis of several considerations. First, it appears to be sufficiently consistent with the remainder to pool. Second, a reanalysis of the pre-1999 studies using the cut point at 4 mG changed the pooled estimate by only 5% (Greenland, Sheppard, Kaune, Poole and Kelsh, 2000), suggesting that the use of 4 rather than 3 mG is of little importance (apart from increasing instability). Third, as will be discussed further, the classifications are at best a surrogate for a true unknown measure, and there are other measurement differences among the studies of potentially much greater importance. Fourth, as with most of the previous studies, covariate adjustment had almost no effect on the estimates.

Two other recent studies were excluded. Green *et al.* (1999) presented only analyses based on quartile categories, resulting in upper cut points of only 1.3-1.5 mG. This study was excluded because the use of such low cut points strongly influenced estimates from earlier studies (Greenland, Sheppard, Kaune, Poole and Kelsh, 2000); it did, however, report positive associations on contrasting the top and bottom quartiles. Schüz *et al.* (2001) had only three highly exposed cases; this study was excluded because of evidence of severe upward bias (twofold or threefold, with odds ratios from 5 to 11) in the reported estimates due to sparse data (Greenland, Schwartzbaum and Finkle, 2000), and because of insufficient reporting of raw data to allow further evaluation.

### 3.2. A conventional analysis

Leukaemia is a very rare disease and the usual justifications for interpreting the observed odds ratios as rate ratio estimates apply (Rothman and Greenland (1998), chapter 7). The odds ratios

are remarkably consistent across studies (homogeneity  $P = 0.24$ ), and the pooled MLE suggests a 70% higher leukaemia rate among children with estimated average exposure above 3 mG. Much like the ML results in Table 1, a Mantel-Haenszel analysis produces an estimated odds ratio for the field-leukaemia association of 1.68, with 95% confidence limits of (1.27, 2.22) and a lower deviance  $P$ -value of 0.0001. Adding study-specific random effects, the usual moment-based overdispersion estimates are 0 owing to the homogeneity, leaving the summary odds ratio and limits virtually unchanged. Under a model with a common rate ratio  $\Omega$  across the underlying study populations, no bias and a uniform prior for  $\theta = \ln(\Omega)$ , the lower  $P$ -value can be interpreted as  $p(\theta < 0|A)$ , the posterior probability that  $\theta < 0$ .

The association is not explained or modified by any known study characteristic or feature of the available data. The results are unchanged by using finer categories (e.g. contrasting greater than 3 mG *versus* less than or equal to 1 mG) or continuous field measurements, and there is no evidence of publication bias (Greenland, Sheppard, Kaune, Poole and Kelsh, 2000). Nonetheless, taking the statistics in Table 1 as unbiased for the field effect is equivalent to assuming that each study reported an experiment in which children were randomized to known residential field levels, were never switched from their initial assignment and were followed until either leukaemia, selection as a control or random censoring occurred.

Put another way, the statistics in Table 1 ignore every source of uncertainty other than random error, including

- (a) possible uncontrolled shared causes (confounders) of field exposure and leukaemia,
- (b) possible uncontrolled associations of exposure and disease with selection and participation (sampling and response biases) and
- (c) magnetic field measurement errors.

Regarding (a), several confounders have been suggested (especially social factors) but there are fewer data on most of these factors than on magnetic fields, and their estimated associations with leukaemia are mostly less than that observed for magnetic fields (to account for the association a factor must by itself have a much stronger association with leukaemia) (Langholz, 2001; Brain *et al.*, 2003). Regarding (b), data suggest that there has been control selection bias in several studies that used direct field measurement (Hatch *et al.*, 2000; Electric Power Research Institute, 2003). Regarding (c), no one doubts that measurement errors must be large. Unfortunately there is no reference measure ('gold standard') for calibration or validation of the measures, in part because no-one knows what an aetiologically relevant magnetic field exposure would be (if one exists). There is only a 'surrogacy' hypothesis that the *known* covariate, contact current, is the 'true' (aetiologically relevant) exposure that is responsible for the observed associations (Kavet and Zaffanella, 2002; Brain *et al.*, 2003), and that magnetic fields are simply an indirect measure of this covariate. Studies are under way to address this hypothesis.

### 3.3. Initial simplifications

Because of the great uncertainty about the bias sources, the inferential situation is very complex and several defensible simplifications will be made. One simplification restricts attention to the dichotomization of field measurements at 3 mG, which greatly eases specification. It was suggested by the repeated observation of almost no field-leukaemia association below 3 mG and was justified by the small changes in conventional statistics that were obtained from a continuous or more finely categorized exposure model (Greenland, Sheppard, Kaune, Poole and Kelsh, 2000; Kabuto, 2003; UK Childhood Cancer Study Investigators, 1999).

Of the three covariates that were uniformly defined and measured on most subjects (study source, age and sex), only study source (modelled as an indicator vector  $S$ ) is used here. On

prior grounds, age and sex among preschool children should be weakly related or unrelated to the exposure and the disease; for example, as noted at least 300 years ago, the sex ratio of births is highly invariant across all factors (Stigler, 1986); also, the rate of leukaemia is only 20% higher among males than females (Brain *et al.*, 2003). Thus, since nearly all the subjects are preschool children, it appears that age and sex can be ignored, and the data conform closely to this expectation; for example, Table 4 of Greenland, Sheppard, Kaune, Poole and Kelsh (2000) shows the small changes in conventional statistics on age–sex adjustment, and adjustment also has little effect in the later studies (Kabuto, 2003; UK Childhood Cancer Study Investigators, 1999). With one exception (London *et al.*, 1991), race is nearly homogeneous within studies and so is automatically controlled by including  $S$  in the models. If further covariate modelling were desired, however, one would expand the vector  $S$  to include other covariates.

Misclassification, non-response and confounding are the bias sources that are believed important by most investigators in this area and will be the only sources that we model. Another source which is often important is publication bias (Copas, 1999), but in the present context such bias is thought to be highly unlikely because of the great public interest in null results (Greenland, Sheppard, Kaune, Poole and Kelsh, 2000). The parameters of the three modelled sources will be given independent priors, so that the full prior covariance matrix is block diagonal with three blocks. This block independence greatly simplifies specification of the prior; misclassification effects on prior information about non-response and confounding will not be considered.

As in almost all the sensitivity analysis literature, confounding will be modelled via a latent variable  $U$  such that the  $US$  conditional field–leukaemia association is unconfounded, i.e. conditioning on  $U$  removes any confounding that is left after conditioning on  $S$  and induces no other confounding. As discussed in Appendix A, the existence of such a sufficient  $U$  is guaranteed under certain causal models, and in special cases this  $U$  may have as few as three levels. Two practical simplifications are made here:  $U$  is further reduced to a binary variable, and the  $US$  conditional field–leukaemia odds ratios are assumed homogeneous across  $U$  given  $S$ . These simplifications greatly ease specification of the prior, do not constrain the amount of confounding in the conventional estimate and appear to have little effect on the results (Greenland, 2003a).

With the above simplifications the classification axes (study variables) are the study, exposure and disease, coded by  $S$ , the row vector of all 14 study indicators,  $X$ , the indicator of field measurement in the top category, and  $Y$ , the leukaemia indicator. The observed data vector  $A$  comprises the  $14(2^2) = 56$   $SXY$  counts of cases and controls in each field category. Define the study level indicators  $D_1 \equiv 1$  if a study used direct (compared with calculated) field measurements and  $V_1 \equiv 1$  if a study was of 120 V–60 Hz (compared with 220 V–50 Hz) systems.  $D_1$  and  $V_1$  also code study location:  $V_1 = 1$  codes North American studies, and  $D_1 = 0$  for all Nordic countries except for the study of Tomenius (1986) (Table 1). Finally, let  $D \equiv (D_1, 1 - D_1)$  and  $V \equiv (V_1, 1 - V_1)$ .  $D$  and  $V$  are functions of  $S$ , so  $S$  contains all the covariate data that are used here.

#### 3.4. Preliminary estimation of uncorrected parameters

As mentioned earlier, certain sequential estimators break down with zero cell counts. Preliminary smoothing eliminates such counts and has theoretical advantages as well (Bishop *et al.*, 1975; Good, 1983). To minimize alteration of data patterns, the observed count vector  $A$  is replaced by a ‘semi-Bayes’ estimate of  $E_\eta$  (Greenland, 2004b), which averages  $A$  (the counts that are expected under a saturated model, regressing  $X$  on  $YS$ ) with those that are expected under a highly saturated model that eliminates 0s, a logistic regression of  $X$  on  $Y, S, YD_1$  and  $YV_1$ . The  $E_\eta$ -estimates that are used here are penalized likelihood fitted values from a mixed effects logistic regression with fixed effects  $Y, S, YD_1$  and  $YV_1$  and normal( $0, \sigma^2$ ) logit residuals

(random effects), where  $\sigma = \ln(20)/2(1.96) = 0.764$ . This estimated  $E_\eta$  is an iterative refinement of averaging the empirical logits (weighted by their inverse empirical variances, with zero weights for undefined logits) with the fixed effects predicted logits (weighted by  $1/\sigma^2$ ); see Greenland (2001b) for a more general example and description of the fitting method. The  $\sigma^2$ -value implies that each exposure odds falls within a 20-fold range of the fixed effects prediction with 95% probability, which is a very weak restriction compared with the fixed-effects-only model. The resulting estimated  $E_\eta$  are very modestly shrunk from  $A$  towards the fixed effects predictions: the largest absolute change in a count is 1.01, the mean absolute change is 0.31, the Mantel–Haenszel statistics are unchanged to the third decimal place and the patterns among study-specific odds ratios (e.g. orderings) are unchanged. This data-structured smoothing should be contrasted with adding a constant to each cell, which is equivalent to averaging observed counts with those fitted from an intercept-only model (Bishop *et al.* (1975), chapter 12).

The uncorrected study-specific odds ratios are now the  $S$ -specific smoothed sample odds ratios

$$\omega_{XY}(s) \equiv E_{11s}E_{00s}/E_{10s}E_{01s}$$

where  $E_{xys}$  is the smoothed count (the estimated  $E_\eta$ -component) at  $X = x$ ,  $Y = y$  and  $S = s$ .  $\hat{\theta}_0$  will be the logarithm of the Mantel–Haenszel weighted average of the  $\omega_{XY}(s)$  over the studies:

$$\omega_{MH0} = \sum_s w_s \omega_{XY}(s) / w_+$$

where  $w_s = E_{10s}E_{01s}/\sum_{xy} E_{xys}$  and  $w_+ = \sum_s w_s$ ;  $\hat{s}_0^2$  will be the standard error estimate of  $\ln(\omega_{MH0})$  of Robins *et al.* (1999) (Rothman and Greenland (1998), page 272). In light of the above discussion regarding efficient weighting, the uncorrected weights  $w_s$  will be used throughout; this fixed weighting over MCSA trials also avoids adding study reweighting effects to bias correction effects in the distribution of corrected estimates.

An alternative that is often used in meta-analysis, the weighted least squares (Woolf) estimator, averages  $\ln\{\omega_{XY}(s)\}$  by using approximate inverse variance weights  $(\sum_{xy} E_{xys}^{-1})^{-1}$  and so (given  $\eta = 0$ ) is first order efficient for the common odds ratio  $\omega$ . In contrast,  $\omega_{MH0}$  is efficient only when  $\omega = 1$ , although it is only slightly inefficient in realistic examples with  $\omega \neq 1$  and exhibits much better behaviour than other estimators when the data are sparse (Breslow, 1981).

### 3.5. Classification error

$X$  will be treated as a misclassified version of a single ‘true’ (but latent) exposure indicator  $T$ . Misclassification correction converts the  $S$ -specific smoothed sample  $XY$  proportions into fitted values for the sample  $TY$  odds ratios  $\omega_{TY}(s)$ ,

$$\omega_{TY}(s) \equiv \frac{p(T=1|Y=1,s)/p(T=0|Y=1,s)}{p(T=1|Y=0,s)/p(T=0|Y=0,s)}, \quad (1)$$

where  $p$  is used to denote sample probability (expected sample proportion). This conversion requires information on the  $TX$ -relationship. Typical prior information concerns the values of the error rates  $p(X=x|T=1-x, y, s)$ . Let  $\varepsilon_0 \equiv p(X=0|T=1, y, s)$  and  $\varepsilon_1 \equiv p(X=1|T=0, y, s)$ , leaving the dependence on  $Y$  and  $S$  implicit; then  $\varepsilon_0$  and  $\varepsilon_1$  are the false negative rates and false positive rates and  $1 - \varepsilon_0 = p(X=1|T=1, y, s)$  and  $1 - \varepsilon_1 = p(X=0|T=0, y, s)$  are the sensitivity and specificity of  $X$  as a measure of  $T$ . Within  $\mathcal{Q}$ , the  $\varepsilon_x$  and  $1 - \varepsilon_x$  are the  $q_{ij}$  within blocks defined by  $S$  and  $Y$ ;  $q_{ij} = 0$  outside those blocks.

It will be assumed that the error rates satisfy the weak condition  $\varepsilon_x < p(X=x|T=x, y, s)$ . The quantity  $p(X=x|y, s)$  is then an identifiable upper bound on  $\varepsilon_x$ , and the low exposure

prevalences that are seen in Table 1 and in surveys imply that the  $\varepsilon_1$  cannot be very large. This does *not* imply that  $X$  is probably correct; for example, we may still have (and often do have)  $p(T = 1|X = 1, y, s) < p(T = 0|X = 1, y, s)$  if the true exposure prevalence  $p(T = 1|y, s)$  is small. All the studies sought to use identical measurement protocols on cases and controls, and in all the studies very high values for  $\varepsilon_x$  (especially  $\varepsilon_1$ ) are implausible. Hence it will be further assumed that within studies the same bound applies to cases and controls, and that this upper bound has a user-specified maximum  $m_x$  across studies:

$$\varepsilon_x < m(x, s) \equiv \min\{p(X = x|Y = 1, s), p(X = x|Y = 0, s), m_x\}. \tag{2}$$

This condition is much weaker than the common assumption that the  $\varepsilon_x$  do not vary with  $Y$  (error ‘non-differential’ with respect to  $Y$ ), although the  $\varepsilon_x$  will be given a very high within-study between- $Y$  correlation. The smoothed data estimates of the  $p(X = x|y, s)$  will be used to estimate the  $m(x, s)$ .

The  $Q$  correction formula (applied to sample expected counts) is equivalent to the standard conversion formula

$$p(T = x|y, s) = \frac{p(X = x|y, s) - \varepsilon_x}{1 - \varepsilon_0 - \varepsilon_1} \tag{3}$$

(Rothman and Greenland (1998), chapter 19), which is positive under the above constraints. The sample  $TY$  odds ratio at  $S = s$  then simplifies to

$$\omega_{TY}(s) = \frac{\{p(X = 1|Y = 1, s) - \varepsilon_1\} / \{p(X = 0|Y = 1, s) - \varepsilon_0\}}{\{p(X = 1|Y = 0, s) - \varepsilon_1\} / \{p(X = 0|Y = 0, s) - \varepsilon_0\}}. \tag{4}$$

Corrections are computed by replacing the  $p(X = x|y, s)$  by the smoothed sample proportions  $E_{1ys} / \sum_x E_{xys}$ , specifying a model for the  $\varepsilon_x$ , and then sampling the model coefficients from a joint prior distribution. At each sampling, the  $\varepsilon_x$  are computed from the model; the corrected sample odds ratios are then computed from the resulting  $\varepsilon_x$ .

There is no quantitative prior information on the error rates. A few studies measured subsets of subjects with different techniques, but the differences in results are highly unstable and there is no evidence on which technique provides a more accurate measure of a true ‘high exposure’ indicator  $T$ . None-the-less, everyone expects considerable heterogeneity in the error rates. Direct and calculated measurements (distinguished by  $D$ ) are vastly different procedures. North American and European power systems (distinguished by  $V$ ) differ in ways that could affect error rates. Because measurement protocols varied greatly across studies, other between-study differences in error rates should also be expected.

The hierarchical misclassification (M) model will thus regress the error rates  $\varepsilon_x$  on  $S$  as well as  $D$  and  $V$ , with  $Y$  included to allow for possible differentiality of errors:

$$\begin{aligned} \eta_M(x|y, s) &\equiv \ln[\varepsilon_x / \{m(x, s) - \varepsilon_x\}] \\ &= \sigma_M \{ \beta_{Mx} + s\beta_{MSx} + d\beta_{MDx} + v\beta_{MVx} + (y, 1 - y)\beta_{MYx} \} \end{aligned} \tag{5}$$

for  $x = 0, 1$ , where  $\eta_M(x|y, s)$  is the logit of the error rate  $\varepsilon_x$  rescaled to the  $\{0, m(x, s)\}$  range of the rate. The model intercepts  $\beta_{Mx}$  have variances  $\tau_{Mx}^2$ . The remaining  $\beta$ -coefficients ( $\beta_{MSx}$ ,  $\beta_{MDx}$ ,  $\beta_{MVx}$ ,  $\beta_{MYx}$ ) represent the dependence of the  $\varepsilon_x$  on the second-stage (group level) covariates  $S$ ,  $D$ ,  $V$  and  $Y$ , and are taken as independent bivariate column vectors whose components are independent with variances ( $\tau^2$ ) such that  $\tau_{Mx}^2 + \tau_{MSx}^2 + \tau_{MDx}^2 + \tau_{MYx}^2 = 1$ . The coefficient scale factor  $\sigma_M$  is a known constant which is introduced solely to make the coefficient variances sum to 1, a feature that eases numerical translation of prior correlations between the  $\eta_M$  into the variance components  $\tau^2$ .

Let  $\beta_M$  denote the vector of all the unknowns (the  $\beta$ ) in the error model, and let  $\omega_{TY}(s; \beta_M)$  be the study-specific corrected odds ratio that is obtained by substituting this model into  $\omega_{TY}(s)$ . To reflect the lack of information for ordering the  $\varepsilon_x$ , I gave all the  $\beta_M$ -components zero means and for convenience gave them normal distributions. To reflect that most of the expected heterogeneity of the  $\varepsilon_x$  is attributed to the type of measurement and study protocol, I assigned  $\tau^2$ -values to produce a simple but plausible prior correlation structure among the  $\eta_M(x|y, s)$ : with  $s_1$  and  $s_0$  coding distinct studies (distinct values of  $S$ ) I wanted correlations of  $\eta_M(x|1, s_1)$  with  $\eta_M(x|0, s_0)$  ranging from small (0.30) among studies that share neither  $D$  nor  $V$  to large (0.70) among studies that share both  $D$  and  $V$ . Within studies, I wanted a nearly perfect case-control correlation (0.95) of  $\eta_M(x|1, s)$  with  $\eta_M(x|0, s)$ ; non-differentiality would correspond to perfect correlation, which is equivalent to dropping  $Y$  from the model. Because there are arguments for positive and negative correlations of the error rates when  $T = 1$  compared with when  $T = 0$ , the  $\eta_M(1|y, s)$  were left uncorrelated with the  $\eta_M(0|y, s)$ ; to induce a correlation we could introduce components that are shared between the coefficients for these two logits. By back-calculation, these choices require  $\tau_{Mx}^2 = \tau_{MVx}^2 = \tau_{MYx}^2 = 0.09$  and  $\tau_{MDx}^2 = 0.31$ , leaving  $\tau_{MSx}^2 = 0.42$ .

Unlike with non-response and confounding, there are no data on which to ground the scale factors and maximum upper bounds  $m_x$  for the  $\varepsilon_x$ . Hence the scale factor  $\sigma_M$  was set to 2, which makes each  $\varepsilon_x$  nearly uniform on its support. The  $m_x$  are the most arbitrary and so will be the focus of a small meta-sensitivity analysis, followed by an analysis in which they are treated as unknown bias parameters.

### 3.6. Non-response

Let  $R(t, y, s)$  be the response rate among population members with  $T = t, Y = y$  and  $S = s$ . The sample probabilities  $p$  are related to the population probabilities  $P$  by

$$p(t_0|y, s) = P(t_0|y, s) R(t_0, y, s) / \sum_t P(t|y, s) R(t, y, s), \tag{6}$$

where the sum is over all possible values of  $T$ . The response bias factors are then

$$B_R(s) \equiv \frac{R(T = 1, Y = 1, s) / R(T = 0, Y = 1, s)}{R(T = 1, Y = 0, s) / R(T = 0, Y = 0, s)}, \tag{7}$$

and hence the population  $TY$  odds ratios

$$\Omega_{TY}(s) \equiv \frac{P(T = 1|Y = 1, s) / P(T = 0|Y = 1, s)}{P(T = 1|Y = 0, s) / P(T = 0|Y = 0, s)} = \frac{P(Y = 1|T = 1, s) / P(Y = 0|T = 1, s)}{P(Y = 1|T = 0, s) / P(Y = 0|T = 0, s)} \tag{8}$$

can be obtained from the sample  $TY$  odds ratios by  $\Omega_{TY}(s) = \omega_{TY}(s) / B_R(s)$ .

Variables that may have important relationships to the response include continent (coded by  $V$ ) and idiosyncrasies of the study design and location (coded by  $S$ ).  $D$  has an expected relationship to the response supported by data on  $X$ : direct measures ( $D_1 = 1$ ) require entry to private property, leading to a low response among controls ( $Y = 0$ ) and among the exposed (Hatch *et al.*, 2000); in contrast, there is high prior probability that studies with calculated fields ( $D_1 = 0$ ) have little or no response bias. Hence the model that is used here is

$$\eta_R(s) \equiv \ln\{B_R(s)\} = \sigma_{RD}(\beta_R + s\beta_{RS} + d\beta_{RD} + v\beta_{RV}) \tag{9}$$

where the  $\beta$ -coefficients represent the dependence of  $B_R$  on the second-stage (group level) covariates  $S, D$  and  $V$ . The variance of the intercept  $\beta_R$  is denoted  $\tau_R^2$ , and the factor coefficients  $\beta_{RS}, \beta_{RD}$  and  $\beta_{RV}$  are taken as independent bivariate column vectors whose components are

independent with variances ( $\tau^2$ ) such that  $\tau_R^2 + \tau_{RS}^2 + \tau_{RD}^2 + \tau_{RV}^2 = 1$ . The scale factor  $\sigma_{RD}$  is treated as known but will depend on  $D$ .

Now let  $\beta_R$  denote the vector of all unknowns ( $\beta$ ) in this specification, and  $B_R(s; \beta_R)$  the response bias model that is obtained by substituting the specification into the bias factor  $B_R(s)$ . To reflect lack of information on response bias sources apart from measurement type, I gave all  $\beta_R$ -components mean 0, except that component 1 of  $\beta_{RD}$  was given mean  $\ln(1.2)/\sigma_{RD}$  on the basis of the elevated non-response among exposed controls that was observed by Hatch *et al.* (2000) and others (Electric Power Research Institute, 2003). For convenience I gave the components normal distributions. To reflect the high prior correlation for non-response across studies, I assigned  $\tau^2$ -values to produce  $\eta_R(s)$  correlations ranging from moderate (0.60) between studies with different  $D$  and  $V$  to very high (0.90) between studies with the same  $D$  and  $V$ ; these are produced by  $\tau_R^2 = \tau_{RD}^2 = 0.36$  and  $\tau_{RV}^2 = 0.09$ , leaving  $\tau_{RS}^2 = 0.19$ . To reflect the greater uncertainty about the amount of response bias in studies with direct measurement, I specified prior 50th, 5th and 95th percentiles for  $B_R(s)$  of 1.20, 0.90 and 1.60 (widely dispersed around 1.2) when  $D_1 = 1$ , and prior 50th, 5th and 95th percentiles for  $B_R(s)$  of 1.00, 0.91 and 1.10 (concentrated around 1) when  $D_1 = 0$ . These percentiles result from assigning  $\sigma_{RD} = \ln(1.33)/1.645$  when  $D_1 = 1$  and  $\sigma_{RD} = \ln(1.10)/1.645$  when  $D_1 = 0$ .

### 3.7. Confounding

Let

$$\Omega_U(s) \equiv P(U = 1|T = Y = 0, s) / P(U = 0|T = Y = 0, s)$$

be the population odds of the latent confounder  $U$  among unexposed non-cases (which compose over 95% of populations in this example), let  $\Omega_{TU}(y, s)$  be the (population)  $TU$  odds ratio given  $YS$ , let  $\Omega_{TY}(u, s)$  be the  $TY$  odds ratio given  $US$  and let  $\Omega_{UY}(t, s)$  be the  $UY$  odds ratio given  $TS$ . As mentioned earlier, I assumed that there is no three-way  $TUY$ -interaction given  $S$ , so that these  $S$ -specific odds ratios are constant over  $T$ ,  $U$  and  $Y$  respectively. The change in the  $S$ -specific  $TY$  odds ratio from ignoring  $U$  is then  $B_C(s) \equiv \Omega_{TY}(s) / \Omega_{TY}(u, s)$ .

Given disease rarity,  $B_C(s)$  is also the degree of  $TY$ -confounding by  $U$  (bias from ignoring  $U$ ) when  $S = s$ . The correction formula is thus  $\Omega_{TY}(u, s) = \Omega_{TY}(s) / B_C(s)$ , where

$$B_C(s) = \frac{\{\Omega_{TU}(y, s) \Omega_{UY}(t, s) \Omega_U(s) + 1\} \{\Omega_U(s) + 1\}}{\{\Omega_{TU}(y, s) \Omega_U(s) + 1\} \{\Omega_{UY}(t, s) \Omega_U(s) + 1\}} \quad (10)$$

(Yanagawa, 1984). By analogy with response bias we could model  $\ln\{B_C(s)\}$  directly (Robins *et al.*, 1999). None-the-less, typical prior information refers instead to the odds ratios in equation (10) and considers those parameters *a priori* independent, which makes it easier to model the  $\Omega$  directly. The models that are used here are

$$\eta_{TU}(s) \equiv \ln\{\Omega_{TU}(s)\} = \sigma_T(\beta_T + s\beta_{TS} + d\beta_{TD} + v\beta_{TV}), \quad (11)$$

$$\eta_U(s) \equiv \ln\{\Omega_U(s)\} = \sigma_U(\beta_U + s\beta_{US} + d\beta_{UD} + v\beta_{UV}), \quad (12)$$

$$\eta_{UY}(s) \equiv \ln\{\Omega_{UY}(s)\} = \sigma_Y(\beta_Y + s\beta_{YS} + d\beta_{YD} + v\beta_{YV}). \quad (13)$$

As with the earlier models, for convenience the linear predictors are rescaled by specified factors  $\sigma_T$ ,  $\sigma_U$  and  $\sigma_Y$  so that the variances ( $\tau^2$ ) of their random ( $\beta$ -) components sum to 1.

Let  $\beta_C$  be the vector of all the  $\beta$  in these three formulae, and let  $B_C(s; \beta_C)$  be the confounding model that is obtained by substituting the specification into the bias factor  $B_C(s)$ . The prior that

is used here is intended to address vague suggestions that some sort of biologically and physically independent leukaemia risk factor may be associated with fields. To reflect the lack of information on specific confounding sources, I gave all  $\beta_C$ -components mean 0 and for convenience gave them normal distributions. Effects of unmeasured factors on leukaemia (parameterized by the  $\eta_{UY}$ ) would be heavily determined by human cancer biology, which is expected to vary little with location, although the distribution of those factors could easily vary. In contrast, the associations of those factors with fields ( $\eta_{TU}$ ) and even more the background prevalences of those factors (whose logits are the  $\eta_U$ ) would be heavily affected by local conditions such as wiring practices. Hence, I assigned  $\tau^2$ -values to produce higher correlations between the  $\eta_{UY}(s)$  than between the  $\eta_{TU}(s)$ , and higher correlations between the  $\eta_{TU}(s)$  than between the  $\eta_U(s)$ . For  $\eta_{UY}(s)$  the correlations ranged from 0.85 between studies that share neither  $D$  nor  $V$  to 0.95 between studies that share both  $D$  and  $V$ , produced by  $\tau_Y^2 = 0.72$  and  $\tau_{YD}^2 = \tau_{YV}^2 = 0.09$ ; for  $\eta_{TU}(s)$  the correlations ranged from 0.60 between studies that share neither  $D$  nor  $V$  to 0.90 between studies that share both  $D$  and  $V$ , produced by  $\tau_T^2 = \tau_{TD}^2 = 0.36$  and  $\tau_{TV}^2 = 0.09$ ; and for  $\eta_U(s)$  the correlations ranged from 0.50 between studies that share neither  $D$  nor  $V$  to 0.70 between studies that share both  $D$  and  $V$ , produced by  $\tau_U^2 = 0.25$  and  $\tau_{UD}^2 = \tau_{UV}^2 = 0.12$ .

The scale factors were based on results of Langholz (2001), who studied 13 factors associated with household wiring in a survey by Bracken *et al.* (1998). Those data exhibited factor prevalences from very low to very high, so  $\sigma_U$  was set to 2 to produce a nearly uniform distribution for  $\Pr(U = 1|x, y, s)$ . The same data also exhibited odds ratios as high as 5.3, which suggests that  $\sigma_T = \ln(6)/1.645$  is reasonable (because this choice makes 6.0 the 95th percentile of the  $\omega_{TU}(s)$  distribution). There are no analogous data on which to base  $\sigma_Y$ , although general observations on the size of composite effects in cancer epidemiology suggest that the symmetrical choice  $\sigma_Y = \sigma_T = \ln(6)/1.645$  is reasonable.

### 3.8. Results from single corrections

The results in Table 2 are based on 250 000 trials for each case and so have Monte Carlo 95% limits within the level of precision that is displayed; hence I shall refer to the observed proportions as probabilities. Before combining corrections, it is instructive to see the effect of each one alone. As a reference point, the first row of Table 2 provides percentiles of the estimated sampling distribution of the uncorrected estimate

$$\omega_{MH0} = \sum_s w_s \omega_{XY}(s) / w_+,$$

which is the distribution of estimates corrected for random error only by drawing a normal(0,  $\hat{\sigma}_0^2$ ) error and subtracting it from  $\ln(\omega_{MH0})$  (Greenland, 2003a). From the 'proportion < 1' column, there is only a 0.01% chance that the random error in the original summary estimate exceeds  $\omega_{MH0}$  (i.e. that random error alone could have moved the summary from less than or equal to 1 to  $\omega_{MH0}$ ).

In an analogous fashion, the second row provides percentiles of estimates  $\sum_s w_s \omega_{XY}(s) / B_R(s; \beta_R)w_+$  corrected for non-response only. Under the above prior, non-response is a much larger source of uncertainty than random error; for example it yields a 5% probability that the net response bias equals or exceeds the observed  $\omega_{MH0}$  (i.e. that non-response alone could have moved the summary from less than or equal to 1 to  $\omega_{MH0}$ ). The third row gives percentiles of the estimates  $\sum_s w_s \omega_{XY}(s) / B_C(s; \beta_C)w_+$  corrected for confounding only. Under the above prior, confounding uncertainty is similar to uncertainty due to random error; for example, there is only a 0.2% probability that the net confounding equals or exceeds  $\omega_{MH0}$  (i.e. that confounding alone could have moved the summary from less than or equal to 1 to  $\omega_{MH0}$ ).



**Table 2.** Percentiles of corrected Mantel–Haenszel odds ratios from multiple-bias analyses of 250 000 trials each, with different maximum upper bounds  $m_0$  and  $m_1$  for the false negative rates  $\varepsilon_0 \equiv p(X=0|T=1, y, s)$  and false positive rates  $\varepsilon_1 \equiv p(X=1|T=0, y, s)$

Factors corrected for	2.5th percentile	50th percentile	97.5th percentile	Proportion <1	Proportion <1.27
Random error only†	1.27	1.68	2.22	0.0001	0.025
Response bias only	0.94	1.45	2.28	0.05	0.29
Confounding only	1.32	1.69	2.33	0.002	0.019
<i>m<sub>0</sub> = 0.05, m<sub>1</sub> = 0.01</i>					
Classification only	1.55	2.07	7.81	<0.0005	<0.0005
All bias sources‡	1.01	1.90	7.32	0.023	0.12
All plus random error§	0.95	1.91	7.50	0.036	0.14
<i>m<sub>0</sub> = 0.05, m<sub>1</sub> = 0.05</i>					
Classification only	1.24	3.63	46.7	0.003	0.031
All bias sources‡	0.96	3.26	42.6	0.031	0.089
All plus random error§	0.92	3.27	43.2	0.037	0.095
<i>m<sub>0</sub> = 0.25, m<sub>1</sub> = 0.01</i>					
Classification only	1.72	2.04	6.95	<0.0005	<0.0005
All bias sources‡	1.06	1.90	6.59	0.015	0.10
All plus random error§	0.99	1.91	6.73	0.027	0.12
<i>m<sub>0</sub> = 0.25, m<sub>1</sub> = 0.05</i>					
Classification only	1.41	3.45	41.3	<0.0005	0.008
All bias sources‡	1.06	3.11	38.1	0.017	0.068
All plus random error§	1.01	3.14	38.6	0.023	0.077
<i>m<sub>0</sub>, m<sub>1</sub> random§§</i>					
Classification only	1.42	2.92	35.1	0.001	0.011
All bias sources‡	1.04	2.67	32.0	0.019	0.077
All plus random error§	0.99	2.70	32.5	0.026	0.088

† Lower 95% limit, point estimate, upper 95% limit and lower *P*-value from a Mantel–Haenszel analysis.  
 ‡ Correcting for bias from misclassification, non-response and confounding.  
 § Adding estimated normal random error (from the first row) to the distribution with all biases.  
 §§  $m_0$  and  $m_1$  logit normal(0, 4) on (0.025, 0.40) and (0.005, 0.105) respectively (roughly uniform on their support).

The first rows of the next four blocks in Table 2 provide percentiles of the estimates  $\Sigma_s w_s \omega_{TY}(s)/w_+$ , corrected for misclassification only, under some reasonable pairs for the maximum bounds  $m_x$  of the  $\varepsilon_x$  across studies. With  $m_0 = 0.05$  and  $m_1 = 0.01$  (which forces  $\varepsilon_0 < 0.05$  and  $\varepsilon_1 < 0.01$  in all studies), the above specification results in a probability of less than 0.05% that the misclassification bias equalled or exceeded  $\omega_{MH0}$  (i.e. that misclassification alone could have shifted the summary estimate from less than or equal to 1 to  $\omega_{MH0}$ ). Increasing  $m_1$  alone to 0.05 increases this probability to 0.3%, but then increasing  $m_0$  to 0.25 reduces the probability to below 0.05% again. In all cases, however, it appears improbable that misclassification alone moved the summary from 1 or less to  $\omega_{MH0}$ .

It may seem paradoxical that increasing classification error bounds can reduce the probability of bias exceeding  $\omega_{MH0}$ . With ‘nearly’ non-differential misclassification, however, the bias that is produced by the misclassification is on average towards 1, in accord with the idea that non-differential exposure measurement error that is independent across units attenuates the

observed association. As a result, the correction to the observed positive association must on average be upwards, and so (as can be seen in Table 2) the corrected estimates are distributed mostly above  $\omega_{MH0}$ , regardless of the bounds. The behaviour of the lower tail of the distribution is more complex, however. First, note that when  $m_0 = m_1 = 0$  there is no misclassification (all  $\varepsilon_x = 0$ ), and so the probability that the classification correction exceeds  $\omega_{MH0}$  is 0. As the  $m_x$  increase from 0 the  $\varepsilon_x$  can vary more widely, and hence the dispersion of the corrected estimates initially expands as the location shifts upwards. The dispersion and location change in a highly non-linear fashion and have opposing effects on the lower tail percentiles. The dispersion can increase more rapidly than the location and thus increase the probability that the correction exceeds  $\omega_{MH0}$  (for example, compare the results for  $m_0 = 0.01$  and  $m_1 = 0.05$  with those for  $m_0 = m_1 = 0.05$ ) but can also decline as the range of misclassification rates increases and thus reduce the probability that the correction exceeds  $\omega_{MH0}$  (for example, compare the results for  $m_0 = 0.25$  and  $m_1 = 0.05$  with those for  $m_0 = m_1 = 0.05$ ). These phenomena can be further explained algebraically but for brevity I omit the details.

### 3.9. Combined corrections, and subsequent inferences

Let  $\beta = (\beta_C, \beta_R, \beta_M)$ . Table 2 provides the percentiles of the multiple-corrected estimates

$$\Omega_{TY}(\beta) = \sum_s w_s \omega_{TY}(s; \beta_M) / B_C(s; \beta_C) B_R(s; \beta_R) w_+$$

for different  $m_x$ -pairs. It also gives percentiles after including log-normal random error at each draw of  $\beta$ , i.e. percentiles of  $\Omega_{TY}(\beta) \exp(Z)$  where  $Z$  is normal( $0, \hat{\sigma}_0^2$ ). We can now look at features of the distribution of the corrected estimates without and with correction for random error and compare these results with those of the conventional analysis (which accounts only for random error). Uncertainty about the *TY*-effect due to uncertainty about classification error is sensitive to the  $m_x$ , especially to the false positive bound  $m_1$  (which is unsurprising, given the low prevalence of exposure). For example, with random error included, the probability that the corrected estimate falls below 1 (i.e. that bias plus random error explain the observed association) is 3.6% for the first pair of  $m_x$  but 2.3% for the last pair. The results are also sensitive to the form of the  $\varepsilon_x$ -distributions within their support (which are not shown). These features should temper any conclusions about the *TY*-effect that might be drawn from the conventional results.

Uncertainty about appropriate bounds suggests adding the  $m_x$  to the model as hyperprior parameters. As an example, the final set of results in Table 2 comes from sampling  $m_0$  and  $m_1$  from logit-normal(0,4) distributions rescaled to (0.025,0.40) and (0.005,0.105) respectively, which are close to uniform on these intervals. The net result of this extension is an averaging of the fixed  $m_x$  results over the range of the  $m_x$  in the sensitivity analysis. Similar results can be obtained by making  $\sigma_M$  unknown with a prior.

Given the prior, the results in Table 2 might be taken as favouring the hypothesis of a leukaemogenic effect of magnetic fields or a close correlate for which they are a surrogate. None-the-less, no agreement about the existence of an effect (let alone its size) could be forced by the data without more precise knowledge of the classification errors. Classification error is the largest source of uncertainty because (unlike non-response and confounding) there are simply no relevant data or theory from which to develop a precise prior for  $\eta_M$ . Even if that information were available, uncertainty that is due to non-response is comparable with uncertainty that is due to random error, and so the overall uncertainty would remain high even if enormous studies with perfect measurements (which will never exist) were added to the analysis. The only positive note is that confounding alone seems to be of lesser importance than other biases, given the prior information that is used here.

Finally, before the earliest studies little credibility was given to the hypothesis that residential fields cause leukaemia. Thus, if we added a substantively justified prior for  $\theta$  to the specification, the final distributions would be shifted towards the null because such a prior is concentrated near the null (Greenland, 2003b).

## 4. Discussion

### 4.1. Model forms

As in most conventional analyses, I have not addressed uncertainty about model forms; I used log-linear and logistic models with normal random effects only for tractability and to enforce range restrictions. This source of uncertainty could be included by adding parameters to index the model form (as is done in Bayesian model averaging), although that would greatly increase the complexity of the bias model and the priors.

In my experience, many users of statistics think that their results do not depend on the form of their model because they use categorized variables or purely tabular analyses. None-the-less, the justification and performance of categorical and tabular statistics depend on implicit regression models; for example, the tabular Cochran–Mantel trend test that is popular in epidemiology is the score test of the slope parameter in a logistic model with binomial errors, and it can be quite misleading when that model form poorly approximates reality (Maclure and Greenland, 1992). Such issues are ordinarily addressed by model diagnostics; on expansion to include bias parameters, however, the model forms (as well as their parameters) are not identified. Thus, in bias modelling the model form is an integral component of the prior specification, rather than a structure that is selected with guidance from the data, as most statistical research treats it.

### 4.2. Some problems in interpretation

Analysts sometimes conclude that the combination of bias and random error is sufficient to explain an elevated estimate if a plausible value or distribution of  $\eta$  could by itself produce a value that is as high as the conventional lower confidence limit; similarly, some analysts call inference about the null sensitive to hidden bias if a plausible value for  $\eta$  could make the two-sided  $P = 0.05$ . These interpretations are misleading because they do not coherently integrate the uncertainties regarding bias and random error. In particular, they suggest that the bias prior makes the null more probable than it actually does. Consider Table 2: the correct probability (under the prior) that the combination of bias and random error equal or exceed the observed association ( $\omega_{MH0}$ ) is the ‘proportion < 1’ in the ‘all plus random error’ row. This probability is always much smaller than the corresponding probability that bias alone could have produced an elevation that is as high or higher than the conventional lower limit (the ‘proportion < 1.27’ in the ‘all bias sources’ row). Given a positive observed association, the use of the lower limit or  $P = 0.05$  as the criterion for evaluating bias sensitivity implicitly assumes that the random error is improbably positive (at least as positive as the difference between the point estimate and the lower limit, an event with only 2.5% probability). These criteria are thus biased in favour of the null hypothesis.

Other analysts report percentiles from MCSA as frequentist statistics; for example, in summarizing the distribution of corrected estimates, they may present the percentage below the null value as a one-sided lower  $P$ -value and refer to the 2.5th and 97.5th percentiles as 95% confidence limits. None-the-less, because the distributions of these summaries have a heavily subjective prior component  $p(\eta)$ , conventional frequentist interpretations are unjustified (Greenland, 2001a).

#### 4.3. Metasensitivity and objections to bias modelling

Bayesian and MCSA outputs depend completely on the prior  $p(\eta)$ , which suggests that a meta-sensitivity analysis of the dependence is essential. Moving in this direction reintroduces the problem of basic sensitivity analysis, however: given the limitless possibilities for  $p(\eta)$ , a thorough metasensitivity analysis would only illustrate how various conclusions can be reached. A conclusion about the target  $\theta$ , however, would require constraints on the  $p(\eta)$ . These constraints would constitute a subjective prior on priors (a hyperprior); incorporating them into the analysis would produce a subjective average of results over the hyperprior, as in the final block of Table 2. This result would itself be subject to concerns about sensitivity to the hyperprior, which would continue on into an infinite regress.

This regress is as unnecessary as it is impractical. Multiple-bias modelling can be treated as a project to discover and exhibit a prior that is arguably reasonable or defensible (in that it is consistent with known facts and established theory), and that leads to borderline conclusive results according to some operative criterion (e.g. a posterior probability for the null of 0.025) (Greenland, 2003a). Such a prior can help to show why the data cannot force agreement between all reasonable observers: defensible perturbations to such a prior can make the results appear moderately inconclusive or moderately conclusive, as the results in Table 2 do when evaluated against a two-sided 0.05-criterion (a criterion that is used in laws and precedents in the USA; see Greenland (2001a)). As an example, in the year following the publication of Greenland, Sheppard, Kaune, Poole and Kelsh (2000) one official of the California State Department of Health publicly asserted, with near certainty, that fields caused childhood leukaemia; this assertion fed demands on the Public Utilities Commission to impose very costly interventions to reduce field levels at schools and homes. Multiple-bias modelling can counter-balance such overconfident assessments of ambiguous evidence and provide more realistic inputs for decision makers (whose decisions will be guided by cost-benefit as well as evidential concerns).

The unlimited nature of metasensitivity may cause some to label bias modelling as a futile exercise. These objections correctly note that, for most topics in which bias modelling might be worthwhile, it would only show how all of an observed association can be plausibly attributed to bias and random error. This objection is no fault of bias modelling, however, but it instead reflects the weakness of available evidence. The demonstration of this weakness is worthwhile if not imperative in many cases, as above.

Metasensitivity has also led to charges that the quantification of uncertainty that is achieved under bias modelling is spurious. There is, however, nothing spurious about the quantification if the prior approximates the views of the analyst, for then the output gives the analyst an idea what his or her posterior bets about the value of  $\theta$  should be. From a more broad perspective, charges of spurious precision embody a double standard relative to the *status quo*: the apparently precise quantification of uncertainty that is offered by conventional analysis is far more spurious than that from bias modelling. Within health sciences, at least, I believe that most researchers fail to grasp how poorly conventional analyses capture uncertainty, and they fail to compensate sufficiently for these deficits. Intuitive discussions of bias often rely on flawed heuristics, such as 'non-differential misclassification introduces a bias toward the null in virtually every study' (Rothman (1986), page 88). Such flawed heuristics ignore the effects of bias uncertainty, effects which are revealed by an exercise in bias modelling (see Section 3.8). A recent sample survey of the epidemiologic literature revealed that most papers do not even deploy flawed heuristics but instead dismiss biases as unlikely to be important, or else simply fail to mention the problems (Jurek *et al.*, 2004). In the rare case that sensitivity analysis is added, it is almost never coherently combined with the assessment of random error.

Another objection is that possible biases are always omitted from modelling. That is true, but the inevitable omission of some bias sources cannot justify the omission of all (which is what conventional analyses do) any more than the inevitable failure to apprehend all murderers can justify ignoring all the murderers who can be apprehended. The inevitability of omissions does suggest that no analysis can do more than to provide a lower bound on the uncertainty that we should have in light of the data and a prior. None-the-less, bias modelling can provide less misleadingly optimistic bounds than can conventional analysis.

4.4. Bias analysis versus better data?

Some recommend that formal bias analysis should be eschewed in favour of improving measurement, response and covariate data. This recommendation is a *non sequitur* and is often wildly impractical. Bias modelling and collection of better data are not mutually exclusive, although bias modelling is often the only feasible option. Exhortations ‘just to collect better data’ are especially empty when (as in the example) we can neither identify a gold standard measurement nor force subjects to participate or to submit to better measurements (which tax co-operation of subjects). Even when we can envision a way to collect better data, decisions must often be made immediately and so can only be based on *currently available* data; as in the example, it may be essential to model those data fully to counter naïve judgments.

‘Collect better data’ becomes a relevant slogan when it is feasible to do so. Multiple-bias modelling is then a useful ally in making clear that the added value of more observations of previous quality (e.g. case-control studies with unknown and possibly large amounts of bias) is much less than conventional statistical formulae convey (Eddy *et al.*, 1992). Conventional standard errors shrink to 0 as the number of observations increases, and total uncertainty approaches the combined bias uncertainty. At some point, mere replication or enlargement of observational studies is not cost effective, and innovations to reduce bias are essential. This point is passed when random error contributes a minority share to total uncertainty. In the example, three more studies of magnetic fields and childhood leukaemia have been published since completion of the pooling project, but none controlled the biases that are described above. Hence, adding these studies has little effect on the final uncertainty distributions; in fact, adding a study with no random error (infinite sample size) but the same bias uncertainty would have little effect.

Most would agree that proposals to confirm or test previous results should include effective safeguards to reduce sources of bias that were suspected in earlier findings, or at least should supply validation data that could provide usefully precise estimates of bias parameters. But the cost of such improvements may be prohibitive. Decisions about funding should also involve considerations of research value (Eddy *et al.*, 1992); the high cost of doing a very informative study may not justify the usual claim that ‘more research is needed’ (Phillips, 2001). When the cost of better data is prohibitive, multiple-bias analysis of existing data may become the best feasible option.

4.5. Concluding remarks

Extreme sensitivity of results to the priors is inevitable and unsurprising, given the many unidentified parameters in realistic bias models. It reflects an irreducible core of uncertainty about the mechanisms generating non-experimental observations (Leamer, 1978; Rubin, 1983). Unfortunately, this core uncertainty is hidden by adherence to identified models. It is more honest instead to bring uncertainty to the fore and to attempt to discover which parameters contribute most to the final uncertainty. Such discovery can help to guide research planning by focusing resources on reducing the largest sources of uncertainty. Those sources are not necessarily the

largest sources of bias, but rather are the sources that are most poorly determined by prior information.

Compared with conventional analysis, multiple-bias modelling better captures uncertainty about effects but requires the specification of a much larger model and demands far more subject-matter knowledge. It also requires much more presentation space and more effort by the reader. Its key advantages may only make it more unappealing: if conducted and presented properly, it depicts how, in the absence of experimental evidence, effects of interest are identified by prior distributions for bias sources rather than by data. It thus belies methods that claim to 'let the data speak for themselves': without external inputs, observational data say nothing at all about causal effects. In many settings it also shows that only indefensibly precise (overconfident) priors can produce firm conclusions, and that conventional methods produce definitive looking results only because they assign probability 1 to the extremely implausible assumption of no bias ( $\eta = 0$ ).

Multiple-bias modelling can be superfluous when conventional standard errors make clear that substantive inferences are unwarranted, as when only a few small studies are available. It may, however, be essential when an analysis purports to draw causal inferences from observational data, when bias uncertainty is comparable with random error or when decisions with costly consequences must be made on the basis of the available evidence (Eddy *et al.*, 1992). I thus argue that bias modelling should become part of the core training of scientists and statisticians who are entrusted with the analysis of observational data. For research planning and allocation, multiple-bias modelling can show when conventional approaches to reducing uncertainty, such as increasing the sample size or replicating studies, have become inefficient (in the magnetic field controversy, this point was reached with studies published in the mid-1990s). To be worthwhile after that point, further studies must give estimates that are more precise and unbiased than previous estimates or else must give precise estimates of biases. When improved studies are prohibitively expensive, multiple-bias modelling may be the best option for decision-making input.

### Acknowledgements

The author is grateful to David Draper, W. Dana Flanders, Patrick Graham, Katherine J. Hoggatt, Tim Lash, George Maldonado, the Joint Editor and the referees for numerous helpful comments, and to Michinori Kabuto for translating results from his study. This research was supported by the Electric Power Research Institute and by grant 1R29-ES07986 from the National Institute of Environmental Health Sciences.

### Appendix A

'Unconfounded' and 'confounding' have been formalized in various ways (Greenland *et al.*, 1999); for the present exposition the precise definition is unimportant as long as it implies that there is a  $U$  such that the true causal effect of  $T$  on  $Y$  can be identified from  $P(T, Y|S, U)$ . This  $U$  may be a compound of other variables. Existence can be shown under various causal models. For example, under a directed acyclic graphical model for causation, all confounding can be traced to common causes of  $T$  and  $Y$ , and hence such a  $U$  will exist if (as here)  $S$  is unaffected by  $T$  or  $Y$  (Pearl (2000), chapter 6). Existence is also guaranteed under a potential outcome model for the effect of  $T$  on  $Y$ , for  $U$  can then be taken as the potential outcome vector (Frangakis and Rubin, 2002). In the present analysis, with binary  $T$ , any sufficient multidimensional  $U$  can be reduced to a sufficient univariate summary; for example, the propensity score  $P(T = 1|S, U)$  is such a summary. This score is usually categorized and five levels are often deemed adequate (Rosenbaum, 2002); if the range of the score is very restricted or the relationship of  $(S, U)$  to  $T$  or  $Y$  is weak, fewer levels may be needed, although more may be needed if the relationship of  $(S, U)$  to  $T$  and  $Y$  is very strong. Note that, under a deterministic monotone effect model for a binary  $Y$ , the potential outcome

vector  $U = (Y_1, Y_0)$  has only three possible levels: (0,0), (1,1) and at most one of (1,0) or (0,1) (Angrist *et al.*, 1996).

## References

- Angrist, J., Imbens, G. and Rubin, D. B. (1996) Identification of causal effects using instrumental variables (with discussion). *J. Am. Statist. Ass.*, **91**, 444–472.
- Bishop, Y. M. M., Fienberg, S. E. and Holland, P. W. (1975) *Discrete Multivariate Analysis: Theory and Practice*. Cambridge: Massachusetts Institute of Technology Press.
- Bracken, M. B., Belanger, K., Hellebrand, K., Adesso, K., Patel, S., Trich, E. and Leaderer, B. (1998) Correlates of residential wiring configurations. *Am. J. Epidemiol.*, **148**, 467–474.
- Brain, J. D., Kavet, R., McCormick, D. L., Poole, C., Silverman, L. B., Smith, T. J., Valberg, P. A., Van Etten, R. A. and Weaver, J. C. (2003) Childhood leukemia: electric and magnetic fields as possible risk factors. *Environ. Hlth Perspect.*, **111**, 962–970.
- Breslow, N. E. (1981) Odds ratio estimators when the data are sparse. *Biometrika*, **68**, 73–84.
- Coghill, R. W., Steward, J. and Philips, A. (1996) Extra low frequency electric and magnetic fields in the bedplace of children diagnosed with leukemia: a case-control study. *Eur. J. Cancer Prev.*, **5**, 153–158.
- Copas, J. B. (1999) What works?: selectivity models and meta-analysis. *J. R. Statist. Soc. A*, **162**, 95–109.
- Copas, J. B. and Li, H. G. (1997) Inference for non-random samples (with discussion). *J. R. Statist. Soc. B*, **59**, 55–95.
- Cornfield, J., Haenszel, W., Hammond, W. C., Lilienfeld, A. M., Shimkin, M. B. and Wynder, E. L. (1959) Smoking and lung cancer: recent evidence and a discussion of some questions. *J. Natn. Cancer Inst.*, **22**, 173–203.
- Crouch, A. C., Lester, R. R., Lash, T. L., Armstrong, S. R. and Green, L. C. (1997) Health risk assessment prepared per the risk assessment reforms under consideration in the U.S. Congress. *Hum. Ecol. Risk Assessmnt*, **3**, 713–785.
- Dockerty, J. D., Elwood, J. M., Skegg, D. C. G. and Herbison, G. P. (1998) Electromagnetic field exposures and childhood cancers in New Zealand. *Cancer Causes Contr.*, **9**, 299–309; erratum **10** (1999), 641.
- Draper, D., Saltelli, A., Tarantola, S. and Prado, P. (2000) Scenario and parametric sensitivity and uncertainty analyses in nuclear waste disposal risk assessment: the case of GESAMAC. In *Mathematical and Statistical Methods for Sensitivity Analysis* (eds A. Saltelli, K. Chan and M. Scott), ch. 13, pp. 275–292. New York: Wiley.
- Eddy, D. M., Hasselblad, V. and Schachter, R. (1992) *Meta-analysis by the Confidence Profile Method*. New York: Academic Press.
- Efron, B. and Tibshirani, R. (1993) *An Introduction to the Bootstrap*. New York: Chapman and Hall.
- Electric Power Research Institute (2003) Selection bias in epidemiologic studies of EMF and childhood leukemia. *EPRI Report 1008149*. World Health Organization, Geneva.
- Feychting, M. and Ahlbom, A. (1993) Magnetic fields and cancer in children residing near Swedish high-voltage power lines. *Am. J. Epidemiol.*, **138**, 467–481.
- Flegal, K. M., Keyl, P. M. and Nieto, F. J. (1991) Differential misclassification arising from nondifferential errors in exposure measurement. *Am. J. Epidemiol.*, **134**, 1233–1244.
- Frangakis, C. and Rubin, D. B. (2002) Principal stratification in causal inference. *Biometrics*, **58**, 21–29.
- Gelman, A., Carlin, J. B., Stern, H. S. and Rubin, D. B. (2003). *Bayesian Data Analysis*, 2nd edn. New York: Chapman and Hall–CRC.
- Good, I. J. (1983) *Good Thinking*. Minneapolis: University of Minnesota Press.
- Graham, P. (2000) Bayesian inference for a generalized population attributable fraction. *Statist. Med.*, **19**, 937–956.
- Green, L. M., Miller, A. B., Villeneuve, P. J., Agnew, D. A., Greenberg, M. L., Li, J. and Donnelly, K. E. (1999) A case-control study of childhood leukemia in southern Ontario, Canada, and exposure to magnetic fields in residences. *Int. J. Cancer*, **82**, 161–170.
- Greenland, S. (1996) Basic methods for sensitivity analysis of bias. *Int. J. Epidemiol.*, **25**, 1107–1116.
- Greenland, S. (1998) The sensitivity of a sensitivity analysis. *Proc. Biometr. Sect. Am. Statist. Ass.*, 19–21.
- Greenland, S. (2001a) Sensitivity analysis, Monte-Carlo risk analysis, and Bayesian uncertainty assessment. *Risk Anal.*, **21**, 579–583.
- Greenland, S. (2001b) Putting background information about relative risks into conjugate priors. *Biometrics*, **57**, 663–670.
- Greenland, S. (2003a) The impact of prior distributions for uncontrolled confounding and response bias: a case study of the relation of wire codes and magnetic fields to childhood leukemia. *J. Am. Statist. Ass.*, **98**, 47–54.
- Greenland, S. (2003b) Generalized conjugate priors for Bayesian analysis of risk and survival regressions. *Biometrics*, **59**, 92–99.
- Greenland, S. (2004a) Interval estimation by simulation as an alternative to and extension of confidence intervals. *Int. J. Epidemiol.*, **33**, 1389–1397.
- Greenland, S. (2004b) Smoothing epidemiologic data. In *Encyclopedia of Biostatistics*, 2nd edn (eds P. Armitage and T. Colton). New York: Wiley.

- Greenland, S., Pearl, J. and Robins, J. M. (1999) Causal diagrams for epidemiologic research. *Epidemiology*, **10**, 37–48.
- Greenland, S. and Robins, J. M. (1985) Confounding and misclassification. *Am. J. Epidemiol.*, **122**, 495–506.
- Greenland, S., Schwartzbaum, J. A. and Finkle, W. D. (2000) Problems from small samples and sparse data in conditional logistic regression analysis. *Am. J. Epidemiol.*, **151**, 531–539.
- Greenland, S., Sheppard, A. R., Kaune, W. T., Poole, C. and Kelsh, M. A. (2000) A pooled analysis of magnetic fields, wire codes, and childhood leukemia. *Epidemiology*, **11**, 624–634.
- Gustafson, P. (2003) *Measurement Error and Misclassification in Statistics and Epidemiology*. New York: Chapman and Hall.
- Hatch, E. E., Kleinerman, R. A., Linet, M. S., Tarone, R. E., Kaune, W. T., Auvinen, A., Baris, D., Robison, L. L. and Wacholder, S. (2000) Do confounding or selection factors of residential wire codes and magnetic fields distort findings of electromagnetic fields studies? *Epidemiology*, **11**, 189–198.
- Jurek, A. M., Maldonado, G., Greenland, S. and Church, T. R. (2004) Exposure-measurement error is frequently ignored when interpreting epidemiologic study results (abstract). *Am. J. Epidemiol.*, **159**, S72.
- Kabuto, M. (2003) A study on environmental EMF and children's health: final report of a grant-in-aid for scientific research project, 1999–2001 (in Japanese). *Report*. Japanese Ministry of Education, Culture, Sports, Science and Technology, Tokyo.
- Kavet, R. and Zaffanella, L. E. (2002) Contact voltage measured in residences: implications for the association between magnetic fields and childhood leukemia. *Bioelectromagnetics*, **23**, 464–474.
- Langholz, B. (2001) Factors that explain the power line configuration wiring code–childhood leukemia association: what would they look like (with discussion)? *Bioelectromagn. Suppl.*, **5**, S19–S31.
- Lash, T. L. and Fink, A. K. (2003) Semi-automated sensitivity analysis to assess systematic errors in observational epidemiologic data. *Epidemiology*, **14**, 451–458.
- Lash, T. L. and Silliman, R. A. (2000) A sensitivity analysis to separate bias due to confounding from bias due to predicting misclassification by a variable that does both. *Epidemiology*, **11**, 544–549.
- Leamer, E. E. (1974) False models and post-data model construction. *J. Am. Statist. Ass.*, **69**, 122–131.
- Leamer, E. E. (1978) *Specification Searches*. New York: Wiley.
- Linet, M. S., Hatch, E. E., Kleinerman, R. A., Robison, L. C., Kaune, W. T., Friedman, D. R., Severson, R. K., Hainer, C. M., Hartsook, C. T., Niwa, S., Wacholder, S. and Tarone, R. E. (1997) Residential exposure to magnetic fields and acute lymphoblastic leukemia in children. *New Engl. J. Med.*, **337**, 1–7.
- Little, R. J. A. and Rubin, D. A. (2002) *Statistical Analysis with Missing Data*, 2nd edn. New York: Wiley.
- London, S. J., Thomas, D. C., Bowman, J. D., Sobel, E., Cheng, T.-C. and Peters, J. M. (1991) Exposure to residential electric and magnetic fields and risk of childhood leukemia. *Am. J. Epidemiol.*, **134**, 923–937.
- Maclure, M. and Greenland, S. (1992) Tests for trend and dose-response: misinterpretations and alternatives. *Am. J. Epidemiol.*, **135**, 96–104.
- Maclure, M. and Schneeweiss, S. (2001) Causation of bias: the episcopo. *Epidemiology*, **12**, 114–122.
- McBride, M. L., Gallagher, R. P., Theriault, H. G., Armstrong, B. G., Tamaro, S., Spinelli, J. J., Deadman, J. E., Fincham, S., Robson, D. and Choi, W. (1999) Power-frequency electric and magnetic fields and risk of childhood cancer. *Am. J. Epidemiol.*, **149**, 831–842.
- Michaelis, J., Schüz, J., Meinert, R., Semann, E., Grigat, J. P., Kaatsch, P., Kaletsh, U., Miesner, A., Brinkmann, K., Kalkner, W. and Kärner, H. (1998) Combined risk estimates for two German population-based case-control studies on residential magnetic fields and childhood leukemia. *Epidemiology*, **9**, 92–94.
- Morgan, M. G. and Henrion, M. (1990) *Uncertainty*. New York: Cambridge University Press.
- Mosteller, F. and Tukey, J. W. (1977) *Data Analysis and Regression*. New York: Addison-Wesley.
- Olsen, J. H., Nielsen, A. and Schulgen, G. (1993) Residence near high voltage facilities and risk of cancer in children. *Br. Med. J.*, **307**, 891–895.
- Pearl, J. (2000) *Causality*. New York: Cambridge University Press.
- Phillips, C. V. (2001) The economics of “more research is needed”. *Int. J. Epidemiol.*, **30**, 771–776.
- Phillips, C. V. (2003) Quantifying and reporting uncertainty from systematic errors. *Epidemiology*, **14**, 459–466.
- Poole, C. and Greenland, S. (1997) How a court accepted a possible explanation. *Am. Statistn*, **51**, 112–114.
- Powell, M., Ebel, E. and Schlossel, W. (2001) Considering uncertainty in comparing the burden of illness due to foodborne microbial pathogens. *Int. J. Food Microbiol.*, **69**, 209–215.
- Robins, J. M., Rotnitzky, A. and Scharfstein, D. O. (1999) Sensitivity analysis for selection bias and unmeasured confounding in missing data and causal inference models. In *Statistical Models in Epidemiology* (eds M. E. Halloran and D. A. Berry), pp. 1–92. New York: Springer.
- Rosenbaum, P. (2002) *Observational Studies*, 2nd edn. New York: Springer.
- Rothman, K. J. (1986) *Modern Epidemiology*. Boston: Little, Brown.
- Rothman, K. J. and Greenland, S. (1998) *Modern Epidemiology*, 2nd edn. Philadelphia: Lippincott.
- Rubin, D. B. (1983) A case study of the robustness of Bayesian methods of inference. In *Scientific Inference, Data Analysis, and Robustness* (eds G. E. P. Box, T. Leonard and C. F. Wu), pp. 213–244. New York: Academic Press.
- Savitz, D. A., Wachtel, H., Barnes, F. A., John, E. M. and Tvrdik, J. G. (1988) Case-control study of childhood cancer and exposure to 60-Hz magnetic fields. *Am. J. Epidemiol.*, **128**, 21–38.



- Schüz, J., Grigat, J. P., Brinkmann, K. and Michaelis, J. (2001) Residential magnetic fields as a risk factor for acute childhood leukemia: results from a German population-based case-control study. *Int. J. Cancer*, **91**, 728–735.
- Steenland, K. and Greenland, S. (2004) Monte-Carlo sensitivity analysis and Bayesian analysis of smoking as an unmeasured confounder in a study of silica and lung cancer. *Am. J. Epidemiol.*, **160**, 384–392.
- Stigler, S. M. (1986) *The History of Statistics*. London: Belknap.
- Tomenius, L. (1986) 50-Hz electromagnetic environment and the incidence of childhood tumors in Stockholm County. *Bioelectromagnetics*, **7**, 191–207.
- Tynes, T. and Haldorsen, T. (1997) Electromagnetic fields and cancer in children residing near Norwegian high-voltage power lines. *Am. J. Epidemiol.*, **145**, 219–226.
- UK Childhood Cancer Study Investigators (1999) Exposure to power-frequency magnetic fields and the risk of childhood cancer. *Lancet*, **354**, 1925–1931.
- Verkasalo, P. K., Pukkala, E., Hongisto, M. Y., Valjus, J. E., Järvinen, P. J., Heikkilä, K. K. and Koskenvuo, M. (1993) Risk of cancer in Finnish children living close to power lines. *Br. Med. J.*, **307**, 895–899.
- Vose, D. (2000) *Risk Analysis*. New York: Wiley.
- Wacholder, S., Armstrong, B. and Hartge, P. (1993) Validation studies using an alloyed gold standard. *Am. J. Epidemiol.*, **137**, 1251–1258.
- Yanagawa, T. (1984) Case-control studies: assessing the effect of a confounding factor. *Biometrika*, **71**, 191–194.

### Discussion on the paper by Greenland

John Copas (*University of Warwick, Coventry*)

It is a pleasure to welcome Professor Greenland to the Society and to propose the vote of thanks for his interesting paper. Most of the work that is reviewed in the paper is published in the epidemiological literature, which is not widely read by statisticians working in other areas. But the central problems of response bias, measurement error and confounding rear their ugly heads in many and probably most applications of statistics, not just in the areas that are usually associated with epidemiology. This is therefore an important topic for all of us.

We are so used to using conventional statistical methods which convert information about a sample  $S$  into a conclusion  $C$  about the population that we all too easily forget the essential impossibility of arguing from the particular to the general. Such induction is only possible if we make assumptions  $A$ , i.e. statistical inference is  $(S, A) \rightarrow C$  and not  $S \rightarrow C$ . Fisher was the first to show that, if we can choose how to obtain  $S$ , then we can do so in such a way that  $A$  is self-evident, in the sense that the randomization that we have actually used to obtain  $S$  also gives us the probability space from which we can obtain  $C$ . If everyone agrees on the truth of  $A$  then there is no need to mention it explicitly. But in all other cases honesty requires us to emphasize that  $C$  depends on  $A$  as well as on  $S$ , and failure to do so is an abuse of our subject. Every day the media invite us to believe claims like ‘studies show that if you eat Corn Flakes for breakfast you are twice as likely to . . .’. No doubt  $S$  is observational, and  $A$  is some absurd assumption of randomization. If  $A$  is not mentioned how can we assess it? If the conclusion is unbelievable, then what is discredited is not  $A$ , as it should be, but statistics, and hence statisticians. I welcome Professor Greenland’s paper for his clear discussion of these issues.

The late George Barnard, in one of these discussion meetings, once remarked ‘We statisticians spend too much time trying to find sensible answers to silly questions’. What question can we ask from Table 2? We see that the upper ends of the intervals for the odds ratio  $\theta$  vary very widely across the different settings for  $(m_0, m_1)$ , but from Professor Greenland’s discussion there seems no very clear reason for preferring any one setting over any other. So, if the question is ‘how big is the risk’, then perhaps we should follow Barnard and answer ‘the quality of the data is not sufficiently good for us to make any sensible estimate’. But, if the question is ‘do we have clear evidence that there is a risk’, we note the remarkable finding that the lower ends of the intervals are all near the null value, even including the setting  $m_0 = m_1 = 0$  if we take this from the ‘response-bias-only’ figures. If we can show that this happens for all reasonable attempts to model these biases, then we have a non-silly question, and the answer is no. This would be an important and perhaps the only convincing analysis of these data.

Professor Greenland argues that sources of bias should be modelled simultaneously and not one by one, as is usually done. Approximating this in terms of sequential transformations on the expected cell counts is an attractive simplification, both conceptually and computationally. So it is a pity that the example does not demonstrate the force of his argument, at least as far as bias is concerned. If we think of each source of bias as adjusting the estimate up or down by a certain factor, and assume that these factors are additive on the log-scale, then we can use the 50% points in the single bias rows of Table 2 to predict the 50% points for the combined bias rows. This gives adjusted values that are very close to the 50% points calculated for the multiple-bias models.